

# Performance Analysis and Optimal Filter Design for Sigma-Delta Modulation via Duality with DPCM

Or Ordentlich and Uri Erez, *Member, IEEE*

**Abstract**—Sampling above the Nyquist rate is at the heart of sigma-delta modulation, where the increase in sampling rate is translated to a reduction in the overall (mean-squared-error) reconstruction distortion. This is attained by using a feedback filter at the encoder, in conjunction with a low-pass filter at the decoder. The goal of this work is to characterize the optimal trade-off between the per-sample quantization rate and the resulting mean-squared-error distortion, under various restrictions on the feedback filter. To this end, we establish a duality relation between the performance of sigma-delta modulation, and that of differential pulse-code modulation when applied to (discrete-time) band-limited inputs. As the optimal trade-off for the latter scheme is fully understood, the full characterization for sigma-delta modulation, as well as the optimal feedback filters, immediately follow.

## I. INTRODUCTION

Analog-to-digital (A/D) and digital-to-analog (D/A) converters are essential in modern electronics. In many cases, it is the quality of these converters that constitutes the main bottleneck in the system, and consequently, dictates its entire performance. On the other hand, as digital circuits are now considered relatively cheap to implement, the interface between the analog and digital domains is often one of the most expensive components in the system. Developing A/D and D/A components that are on the one hand relatively simple, and on the other hand introduce little distortion, is therefore of interest.

Often, the same A/D (or D/A) component is applied to a variety of signals with distinct characterizations. For this reason, it is desirable to design the data converter to be robust to the characteristics of the input signal. One assumption that cannot be avoided is the bandwidth of the signal to be converted, which dictates the minimal sampling rate, according to Nyquist's theorem. Beyond bandwidth, however, one would like to assume as little as possible about the input signal. A reasonable model for the input signal is therefore a *stochastic* one, where the input signal is assumed to be a stationary Gaussian process with a given variance and an arbitrary *unknown* power spectral density (PSD) within the assumed bandwidth, and zero otherwise. In this paper, we adopt this *compound* model which is rich enough to include a wide variety of processes. The robustness requirement from the A/D

(or D/A) converter translates to requiring that it induces a small average distortion simultaneously for all processes within our compound model.

Sigma-delta modulation is a widely used technique for A/D as well as D/A conversion. The main advantage offered by this type of modulation is the ability to trade-off the sampling rate and the number of bits per sample required for achieving a target mean-squared error (MSE) distortion. The input to the sigma-delta modulator is a signal sampled at  $L$  times the Nyquist rate ( $L > 1$ ). This over-sampled signal is then quantized using an  $R$ -bit quantizer. In much of the literature about sigma-delta modulation, no stochastic model is assumed for the input signal. However, when such a model is assumed, the benefit of over-sampling can be easily understood from basic rate-distortion theoretic principles: the (per-sample) rate required to achieve distortion  $D$  for the over-sampled signal is  $L$  times smaller than the rate required to achieve the same distortion for the signal obtained by sampling at the Nyquist rate. Thus, in principle, increasing the sampling rate should allow one to use quantizers with lower resolution, which is desirable in many applications.

However, the rate-distortion theoretical property that guarantees a constant product of the number of bits per sample needed to achieve distortion  $D$ , and the over-sampling ratio  $L$ , is only valid when a very long block of samples is vector-quantized. In A/D and D/A conversion, vector-quantization in high dimensions is a prohibitively complex operation, and quantization is invariably done via scalar uniform quantizers. Scalar quantizers alone cannot translate the increase of sampling rate to a significant reduction in the necessary resolution, but fortunately this problem can be circumvented with the aid of appropriate signal processing.

In sigma-delta based converters, the quantization noise is shaped using a causal shaping filter embedded within a feedback loop, see Figure 1. The filter coefficients are chosen in a manner that ensures that most of the energy of the shaped quantization noise lies outside the frequency band occupied by the over-sampled signal. At the decoder, the quantized signal is low-pass filtered, cancelling out the high-frequencies of the quantization noise process without effecting the signal, such that the decoder's output is composed of the original signal corrupted by a low-pass noise process.

Another technique for compressing sources with memory, which explicitly models the source as a stochastic process, is differential pulse-code modulation (DPCM). In DPCM, a prediction filter is applied to the quantized signal. The output of this filter is then subtracted from the source and the result is fed to the quantizer, see Figure 2. At the decoder, the

The work of O. Ordentlich was supported by the Adams Fellowship Program of the Israel Academy of Sciences and Humanities, a fellowship from The Yitzhak and Chaya Weinstein Research Institute for Signal Processing at Tel Aviv University and the Feder Family Award. The work of U. Erez was supported by the ISF under Grant 1557/13.

O. Ordentlich and U. Erez are with Tel Aviv University, Tel Aviv, Israel (email: ordent,uri@eng.tau.ac.il).

quantized signal is simply passed through the inverse of the prediction filter. The well-known “DPCM error identity” [1] states that the output of the decoder is equal to the source plus the quantization error, just like in simple non-predictive quantization. The benefit of using DPCM, however, is that the signal fed to the quantizer is the error in predicting the source from its *quantized* past, rather than the source itself. If the coefficients of the prediction filter are chosen appropriately, the variance of this error should be smaller than the variance of the original source, which translates to a reduction in the number of bits required from the quantizer for achieving a certain distortion.

The performance of DPCM under the assumption of high-resolution quantization is well understood since as early as the mid 60’s [1]–[3]. Under this assumption, the prediction filter should be chosen as the optimal linear minimum mean-squared-error (MMSE) prediction filter of the source process from its past [1], and the effect of the filtered quantization noise can be neglected in the prediction process. While in most cases where DPCM is traditionally used, the high resolution assumption is well justified, it totally breaks down for the class of band-limited processes, which includes the input signals to sigma-delta modulators. Indeed, the prediction error of such a process from its infinite past has zero-variance, rendering the DPCM high-resolution rate-distortion formulas completely useless.

#### A. Connection to Previous Work

The connection between DPCM and sigma-delta modulation, as two instances of predictive coding, was known from the outset. Indeed, both paradigms emerged from two Bell-Labs patents authored by CC Cutler [4], [5] in 1952 and 1954.

In fact, by adding appropriate pre- and post-filters to the sigma-delta modulator, as depicted in Figure 3, the input to the quantizer, as well as the final reconstruction of the signal, become identical to those in the DPCM architecture [6, Section II], [7, Chapter 3.2.4]. For this reason, it has become folklore that the two architectures are equivalent. When a sigma-delta modulator is used for compression of digital discrete-time signals, the pre-filtering can be performed digitally and the additional complexity of the architecture depicted in Figure 3, w.r.t. that in depicted in Figure 1, may be acceptable. This is however *not* the case for data converters, as the input to the latter is analog and pre-filtering must be done in continuous-time, which is more challenging. The motivation for this work is understanding the performance limits of A/D and D/A conversion based on the sigma-delta architecture, and therefore pre-filtering is precluded. Thus, the architecture is confined to that depicted in Figure 1.

Another important aspect of our interest in sigma-delta modulators as a mean of data-conversion rather than data-compression, is that it dictates that the assumptions one can make on the statistics of the input signal must be minimal. Consequently, we consider a *compound* class of sources that consists of all stationary Gaussian processes with variance  $\sigma_X^2$  whose PSD is limited to some predefined frequency band. In

addition, since data converters often operate at very high rates, it makes sense to impose various constraints on the sigma-delta feedback filter  $C(Z)$ , such as confining it to be a finite impulse response (FIR) filter with a limited number of taps. For a given desired MSE distortion level, our goal is to find the constrained sigma-delta feedback filter  $C(Z)$  that minimizes the quantization rate w.r.t. all sources in the compound class, and to characterize the attained rate. This goal is different than the one pursued in [8], where the optimal *unconstrained* filters w.r.t. a known PSD were found.

The problem of finding the optimal  $N$ -tap FIR sigma-delta feedback filter  $C(Z)$  for a compound family of sources similar to ours, was considered in [6]. The optimal filter was claimed in [6] to be the  $N$ th order MMSE prediction filter  $C(Z) = (1 - Z^{-1})^N$  of a bandpass stationary process from its past, and for a fixed target MSE distortion the required quantization rate was found to decrease linearly with  $N$ . Such a statement is obviously inaccurate, as it violates Shannon’s rate-distortion theorem. The major drawback of [6] is that it (implicitly) makes the high-resolution assumption that the variance of the quantizer’s input is solely dictated by the target signal  $\{X_n\}$ , whereas the contribution of the quantization noise to this variance can be neglected. As discussed above, for over-sampled processes this assumption may not be valid even when the quantizer’s resolution is very high. In particular, using the filter  $C(Z) = (1 - Z^{-1})^N$  from [6], the energy of the quantization noise within the frequency band occupied by the signal indeed decreases exponentially with  $N$ . However, the noise’s energy outside this band increases rapidly with  $N$ , and for any quantization resolution it will become much greater than  $\sigma_X^2$  for  $N$  large enough, making the high-resolution assumption inapplicable. In this case, the dynamic range of the quantizer will be exceeded and overload errors would frequently occur.

It therefore follows that in the analysis of sigma-delta modulators one should not make high-resolution assumptions, but rather must take into account the effect of the filtered quantization noise on the variance of the quantizer’s input. Fortunately, in the analysis of DPCM modulators the high-resolution assumption has been overcome in [9]. It was shown that for any distortion level and any stationary Gaussian source, the DPCM architecture induces a rate-distortion optimal test channel, provided that the prediction filter is chosen as the optimal filter for predicting the source from its *quantized past*, and in addition water-filling pre- and post-filters are applied. The analysis of [9], which takes into account the effect of the quantization noise, can therefore be used to obtain the optimal feedback filter and its corresponding performance for a DPCM system applied to an over-sampled stationary Gaussian source. In this paper, we leverage the results from [9] to the analysis of sigma-delta modulators, by establishing an appropriate duality between the two architectures.

#### B. Contributions

Let  $\mathcal{S}$  be the compound class of all discrete-time stationary Gaussian sources with variance  $\sigma_X^2$  and PSD that is zero for

all  $\omega \notin [-\pi/L, \pi/L]$ ,  $L \geq 1$ . Note that this class corresponds to uniformly sampling a compound class of continuous-time stationary Gaussian processes with variance  $\sigma_X^2$  and PSD that is zero for all  $|f| > f_{\max}$ , at a sampling rate of  $2Lf_{\max}$  samples/per second. Let  $\{X_n^{\text{DPCM}}\}$  be a discrete-time stationary Gaussian process with PSD

$$S_X^{\text{DPCM}}(\omega) = \begin{cases} L\sigma_X^2 & \text{for } |\omega| \leq \pi/L \\ 0 & \text{for } \pi/L < |\omega| < \pi \end{cases}, \quad (1)$$

and note that  $\{X_n^{\text{DPCM}}\} \in \mathcal{S}$ .

Our main result, derived in Section II, is that for any process  $\{X_n^{\Sigma\Delta}\}$  from the compound class  $\mathcal{S}$ , the test channel induced by the sigma-delta modulator (Figure 1) achieves exactly the same rate-distortion function as that of the DPCM test channel (Figure 2) with input  $\{X_n^{\text{DPCM}}\}$ . More specifically, for such processes, for any choice of  $\sigma_{\text{DPCM}}^2$  and prediction filter  $C(Z)$  in the test channel of Figure 2, the same choice of  $C(Z)$  together with the choice

$$\sigma_{\Sigma\Delta}^2 = \frac{\sigma_{\text{DPCM}}^2}{L \cdot \frac{1}{2\pi} \int_{-\pi/L}^{\pi/L} |1 - C(\omega)|^2 d\omega} \quad (2)$$

in Figure 1, yields the same compression rate and the same distortion.

While this result is simple to derive, it has a very pleasing consequence: the problem of optimizing the filter  $C(Z)$  in sigma-delta modulation w.r.t. any signal in  $\mathcal{S}$ , under any set of constraints, can be cast as an equivalent problem of optimizing the DPCM prediction filter w.r.t. input  $\{X_n^{\text{DPCM}}\}$  under the same set of constraints. Furthermore, in Section II-A, we formalize a similar duality between DPCM and sigma-delta modulation for a frequency-weighted-mean-squared-error distortion measure. In this case  $S_X^{\text{DPCM}}(\omega)$  is replaced with a PSD that depends on the distortion's weight function.

In principle, recasting the sigma-delta optimization problem as an MMSE prediction problem may be derived directly from the formulas characterizing its performance, as given in Proposition 1. Nevertheless, establishing the equivalence between sigma-delta modulation and DPCM, in the specific form described above, is insightful as it allows to borrow known results from the literature about the latter.

Having recast the filter optimization problem for sigma-delta as that of optimal linear prediction, we can readily obtain the solution under constraints for which an explicit solution was lacking in the literature, or was cumbersome to derive.

One may question the relevance of the test channel of Figure 1 and its information-theoretic analysis to the practical, resource limited, problem of A/D and D/A conversion via sigma-delta modulators. To that end, in Section III we replace the AWGN channel from Figure 1 with a simple scalar uniform (dithered) quantizer of finite support, which is suitable for implementation within A/D and D/A converters. As long as overload does not occur, the effect of applying the scalar quantizer is equivalent to that of an additive noise channel. We show that the rate-distortion trade-off derived for sigma-delta modulation in Section II remains valid with high probability, with a constant additive excess-rate penalty for using scalar

quantization. The purpose of this excess-rate is to ensure that an overload event, which jeopardizes the stability of the system, occurs with low probability. The stochastic model we assume for the input process allows us to tackle the issue of stability in a systematic and rigorous manner, and the trade-off between the excess-rate and the overload probability is analytically determined.

Clearly, a sigma-delta modulator can only perform well if overload errors are rather rare. Our stability analysis in Section III is based on avoiding overload events w.h.p., and does not aim to consider the effect of such events on the distortion once they occur. In general, the overload probability of the scheme described in Section III decreases double exponentially with the excess-rate of the quantizer w.r.t. the mutual information. Thus, taking an excess rate of 1 – 2 bits will usually yield a sufficiently low overload probability. However, sigma-delta quantizers are often employed with a one-bit quantizer. In this case, the overload error probability cannot be very low. Consequently, the designer would need to guarantee that the effect of overload errors is local in time, and does not drive the system out of stability. There are various restrictions one can place on  $C(Z)$  in pursuit of the latter goal. The issue of maintaining stability when overload errors are unavoidable is outside the scope of this paper. Nevertheless, we stress that our main result is of great relevance to this setting, as it shows that the filter  $C(Z)$  should be chosen as the optimal MMSE prediction filter of  $\{X_n^{\text{DPCM}}\}$  from its noisy past under the stability ensuring restrictions.

## II. MAIN RESULT

We begin by introducing some basic notation that will be used in the sequel. For a discrete signal  $\{c_n\}$ , the  $Z$ -transform is defined as

$$C(Z) \triangleq \sum_{n=-\infty}^{\infty} c_n Z^{-n},$$

and the Fourier transform as

$$C(\omega) \triangleq C(Z)|_{Z=e^{j\omega}} = \sum_{n=-\infty}^{\infty} c_n e^{-j\omega n}.$$

For a discrete (real) stationary process  $\{X_n\}$  with zero-mean and autocorrelation function  $R_X[k] \triangleq \mathbb{E}(X_{n+k}X_n)$  we define the power-spectral density (PSD) as the Fourier transform of the autocorrelation function

$$S_X(\omega) \triangleq \sum_{k=-\infty}^{\infty} R_X[k] e^{-j\omega k}.$$

The PSD of a continuous stationary process is defined in an analogous manner.

Assume  $X^{\Sigma\Delta}(t)$  is a continuous stationary band-limited Gaussian process with zero mean and variance  $\sigma_X^2$ , whose PSD is zero for all frequencies  $|f| > f_{\max}$ , but otherwise unknown. The Nyquist sampling rate for this process is  $2f_{\max}$  samples per second. Since our focus here is on quantization of over-sampled signals, we assume that  $X^{\Sigma\Delta}(t)$  is sampled uniformly with rate of  $2Lf_{\max}$  samples per second for some  $L > 1$ .



quantizer of finite support.

We begin with the test channel in Figure 1, that corresponds to sigma-delta modulation, with the sigma-delta quantizer replaced by an AWGN channel with zero mean and variance  $\sigma_{\Sigma\Delta}^2$ . The filter  $C(Z)$  is assumed to be strictly causal.

*Proposition 1:* For any Gaussian stationary process  $\{X_n^{\Sigma\Delta}\}$  with variance  $\sigma_X^2$  whose PSD is zero for all  $\omega \notin [-\pi/L, \pi/L]$ , the test channel from Figure 1 achieves MSE distortion

$$D = \sigma_{\Sigma\Delta}^2 \cdot \frac{1}{2\pi} \int_{-\pi/L}^{\pi/L} |1 - C(\omega)|^2 d\omega,$$

and its scalar mutual information satisfies<sup>1</sup>

$$\begin{aligned} I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) \\ = \frac{1}{2} \log \left( 1 + \frac{1}{2\pi} \int_{-\pi}^{\pi} |C(\omega)|^2 d\omega + \frac{\sigma_X^2}{\sigma_{\Sigma\Delta}^2} \right). \end{aligned}$$

*Proof:* From Figure 1, we have that

$$U_n^{\Sigma\Delta} = X_n^{\Sigma\Delta} - c_n * N_n^{\Sigma\Delta}, \quad (3)$$

and therefore

$$U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta} = X_n^{\Sigma\Delta} + (\delta_n - c_n) * N_n^{\Sigma\Delta},$$

where  $\delta_n$  is the discrete identity filter. Using the fact that  $\{X_n^{\Sigma\Delta}\}$  is a low-pass process, passing it through the filter  $H(\omega)$  has no effect, and hence

$$\begin{aligned} \hat{X}_n^{\Sigma\Delta} &= h_n * (U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) \\ &= X_n^{\Sigma\Delta} + h_n * (\delta_n - c_n) * N_n^{\Sigma\Delta}. \end{aligned}$$

The MSE distortion attained by the test channel from Figure 1 is therefore

$$D = \mathbb{E}(X_n^{\Sigma\Delta} - \hat{X}_n^{\Sigma\Delta})^2 = \sigma_{\Sigma\Delta}^2 \cdot \frac{1}{2\pi} \int_{-\pi/L}^{\pi/L} |1 - C(\omega)|^2 d\omega.$$

The scalar mutual information between the “quantizer’s” input  $U_n^{\Sigma\Delta}$  and output  $U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}$  is given by

$$I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) = h(U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) - h(N_n^{\Sigma\Delta}) \quad (4)$$

$$= \frac{1}{2} \log \left( 1 + \frac{\mathbb{E}(U_n^{\Sigma\Delta})^2}{\sigma_{\Sigma\Delta}^2} \right), \quad (5)$$

where (4), as well as (5), follow from the statistical independence of  $N_n^{\Sigma\Delta}$  and  $U_n^{\Sigma\Delta}$ . Using (3), the variance of  $U_n^{\Sigma\Delta}$  is

$$\mathbb{E}(U_n^{\Sigma\Delta})^2 = \sigma_X^2 + \sigma_{\Sigma\Delta}^2 \frac{1}{2\pi} \int_{-\pi}^{\pi} |C(\omega)|^2 d\omega. \quad (6)$$

Substituting (6) into (5) establishes the second part of the proposition. ■

Next, we analyze the test channel in Figure 2, that corresponds to DPCM compression with the DPCM quantizer replaced by an AWGN channel with zero mean and variance  $\sigma_{\text{DPCM}}^2$ . As in the test channel of Figure 1, the filter  $C(Z)$  is strictly causal. The distortion corresponding to this test channel, as well as  $I(U_n^{\text{DPCM}}; U_n^{\text{DPCM}} + N_n^{\text{DPCM}})$ , were already found in [9, Theorem 1] for the special case where  $C(Z)$  is

the optimal MMSE infinite length prediction filter of  $X_n^{\text{DPCM}}$  from all past samples of the process  $\{X_n^{\text{DPCM}} + N_n^{\text{DPCM}}\}$ . The following straightforward proposition characterizes the rate and distortion for any choice of the causal filter  $C(Z)$  and any value of  $\sigma_{\text{DPCM}}^2$ .

*Proposition 2:* For a Gaussian stationary process  $\{X_n^{\text{DPCM}}\}$  with variance  $\sigma_X^2$  and PSD

$$S_X^{\text{DPCM}}(\omega) = \begin{cases} L\sigma_X^2 & \text{for } |\omega| \leq \pi/L \\ 0 & \text{for } \pi/L < |\omega| < \pi \end{cases}, \quad (7)$$

the test channel from Figure 2 achieves MSE distortion

$$D = \frac{\sigma_{\text{DPCM}}^2}{L},$$

and its scalar mutual information satisfies

$$\begin{aligned} I(U_n^{\text{DPCM}}; U_n^{\text{DPCM}} + N_n^{\text{DPCM}}) &= \frac{1}{2} \log \left( 1 + \frac{1}{2\pi} \int_{-\pi}^{\pi} |C(\omega)|^2 d\omega \right. \\ &\quad \left. + \frac{L\sigma_X^2}{\sigma_{\text{DPCM}}^2} \frac{1}{2\pi} \int_{-\pi/L}^{\pi/L} |1 - C(\omega)|^2 d\omega \right). \end{aligned}$$

*Proof:* From Figure 2, we have that

$$U_n^{\text{DPCM}} = X_n^{\text{DPCM}} - c_n * V_n^{\text{DPCM}} \quad (8)$$

$$V_n^{\text{DPCM}} = U_n^{\text{DPCM}} + N_n^{\text{DPCM}} + c_n * V_n^{\text{DPCM}} \quad (9)$$

Substituting (8) in (9) yields

$$V_n^{\text{DPCM}} = X_n^{\text{DPCM}} + N_n^{\text{DPCM}}. \quad (10)$$

Using the fact that  $\{X_n^{\text{DPCM}}\}$  is a low-pass process, as before, we obtain

$$\begin{aligned} \hat{X}_n^{\text{DPCM}} &= h_n * (X_n^{\text{DPCM}} + N_n^{\text{DPCM}}) \\ &= X_n^{\text{DPCM}} + h_n * N_n^{\text{DPCM}}. \end{aligned} \quad (11)$$

Since  $\{N_n^{\text{DPCM}}\}$  is AWGN with variance  $\sigma_{\text{DPCM}}^2$ , the variance of the filtered process  $h_n * N_n^{\text{DPCM}}$  is  $\sigma_{\text{DPCM}}^2/L$ . Thus,

$$D = \mathbb{E}(X_n^{\text{DPCM}} - \hat{X}_n^{\text{DPCM}})^2 = \frac{\sigma_{\text{DPCM}}^2}{L}.$$

As in the analysis of the test channel from Figure 1, the scalar mutual information between  $U_n^{\text{DPCM}}$  and  $U_n^{\text{DPCM}} + N_n^{\text{DPCM}}$  is given by

$$I(U_n^{\text{DPCM}}; U_n^{\text{DPCM}} + N_n^{\text{DPCM}}) = \frac{1}{2} \log \left( 1 + \frac{\mathbb{E}(U_n^{\text{DPCM}})^2}{\sigma_{\text{DPCM}}^2} \right). \quad (12)$$

Now, substituting (10) in (8) gives

$$U_n^{\text{DPCM}} = (\delta_n - c_n) * X_n^{\text{DPCM}} - c_n * N_n^{\text{DPCM}},$$

and the variance of  $U_n$  is therefore

$$\begin{aligned} \mathbb{E}(U_n^{\text{DPCM}})^2 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} S_X^{\text{DPCM}}(\omega) |1 - C(\omega)|^2 d\omega \\ &\quad + \frac{1}{2\pi} \int_{-\pi}^{\pi} S_N^{\text{DPCM}}(\omega) |C(\omega)|^2 d\omega \\ &= \frac{L\sigma_X^2}{2\pi} \int_{-\pi/L}^{\pi/L} |1 - C(\omega)|^2 d\omega + \frac{\sigma_{\text{DPCM}}^2}{2\pi} \int_{-\pi}^{\pi} |C(\omega)|^2 d\omega. \end{aligned} \quad (13)$$

<sup>1</sup>All logarithms in this paper are taken to base 2.

Substituting (13) into (12) establishes the second part of the proposition. ■

*Remark 1:* In propositions 1 and 2 we derived the scalar mutual information between the input and output of the AWGN test channels embedded in Figures 1 and 2, respectively. As will become clear in Section III, the scalar mutual information is closely related to the required quantization rate when a scalar memoryless quantizer is used within the sigma-delta or DPCM modulator. In [9], [11], the directed information was shown to be related to the required quantization rate when the quantizer is followed by an entropy coder. Here, we do not consider applying entropy coding to the quantizer's output as we require that the designed modulator be robust to the statistics of the input process, whereas entropy coding is very sensitive to the process statistics. Moreover, if the design of an A/D (or D/A) is considered, the appropriate merit for the modulator's complexity is the number of quantization levels within the scalar quantizer, which are not reduced by incorporating an entropy coder.

Our main result now follows immediately from Propositions 1 and 2.

*Theorem 1:* Let  $\{X_n^{\Sigma\Delta}\}$  be any Gaussian stationary process with variance  $\sigma_X^2$  whose PSD is zero for all  $\omega \notin [-\pi/L, \pi/L]$ , let  $\{X_n^{\text{DPCM}}\}$  be a flat low-pass Gaussian stationary process with PSD as in (7), and let  $C(Z)$  be a strictly causal filter. The test channel from Figure 1 with

$$\sigma_{\Sigma\Delta}^2 = \frac{D}{\frac{1}{2\pi} \int_{-\pi/L}^{\pi/L} |1 - C(\omega)|^2 d\omega},$$

and the test channel from Figure 2 with

$$\sigma_{\text{DPCM}}^2 = L \cdot D,$$

both achieve MSE distortion  $D$  and their scalar mutual information satisfy

$$\begin{aligned} I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) &= I(U_n^{\text{DPCM}}; U_n^{\text{DPCM}} + N_n^{\text{DPCM}}) \\ &= \frac{1}{2} \log \left( 1 + \frac{1}{2\pi} \int_{-\pi}^{\pi} |C(\omega)|^2 d\omega \right. \\ &\quad \left. + \frac{\sigma_X^2}{D} \frac{1}{2\pi} \int_{-\pi/L}^{\pi/L} |1 - C(\omega)|^2 d\omega \right). \end{aligned}$$

This theorem indicates that for any stationary band-limited Gaussian process with variance  $\sigma_X^2$ , the sigma-delta test channel from Figure 1 achieves exactly the same rate-distortion trade-off as that of the DPCM test channel from Figure 2 with a stationary flat band-limited Gaussian input with the same variance, provided that the AWGN variances are scaled according to (2). Thus, Theorem 1 provides a unified framework for analyzing the performance of sigma-delta modulation and DPCM. A great advantage offered by such a unified framework, is that any result known for DPCM can be translated to a corresponding result for sigma-delta modulation, and vice versa. Theorems 2 and Corollary 1 below constitute two important examples of such results.

*Theorem 2:* Let  $\{X_n^{\Sigma\Delta}\}$  be a Gaussian stationary process with variance  $\sigma_X^2$  whose PSD is zero for all  $\omega \notin [-\pi/L, \pi/L]$  and let  $\mathcal{C}$  be a family of strictly causal filters. Define the “virtual” process  $\{S_n\}$  as a Gaussian stationary process with PSD as in (7), and the “virtual” process  $\{W_n\}$  as a Gaussian i.i.d. random process statistically independent of  $\{S_n\}$  with variance  $L \cdot D$ ,  $D > 0$ . Let

$$\begin{aligned} \sigma_D^{*2} &= \min_{C(Z) \in \mathcal{C}} \mathbb{E} (S_n - c_n * (S_n + W_n))^2 \\ C_D^*(Z) &= \operatorname{argmin}_{C(Z) \in \mathcal{C}} \mathbb{E} (S_n - c_n * (S_n + W_n))^2. \end{aligned}$$

If the filter  $C(Z)$  in the sigma-delta test channel from Figure 1 belongs to  $\mathcal{C}$  and the MSE distortion attained by this test channel is  $D$ , then

$$I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) \geq \frac{1}{2} \log \left( 1 + \frac{\sigma_D^{*2}}{L \cdot D} \right), \quad (14)$$

with equality if  $C(Z) = C_D^*(Z)$ .

Theorem 2 states that for a target distortion  $D$ , the sigma-delta filter which minimizes the required compression rate is the optimal linear time-invariant MMSE estimator, within the class of constraints  $\mathcal{C}$ , for  $S_n$  from the past of the noisy process  $\{S_n + W_n\}$ . For example, if  $\mathcal{C}$  consists of all strictly causal finite-impulse response (FIR) filters of length  $p$ , the optimal filter  $C(Z)$  is the optimal predictor of  $S_n$  from the samples  $\{S_{n-1} + W_{n-1}, \dots, S_{n-p} + W_{n-p}\}$ , which can be easily calculated in closed-form.

The optimal sigma-delta filter design problem was studied by several authors, under various assumptions [1], [6], [8], [12]–[15]. However, to the best of our knowledge, the simple expression from Theorem 2 for the optimal filter as the optimal predictor of  $S_n$  from the past of  $\{S_n + W_n\}$  is novel. The references most relevant to Theorem 2, are perhaps [14] and [8], [15]. In [14], Spang and Schultheiss formulated an optimization problem for finding the best FIR filter with  $p$  coefficients in a sigma-delta modulator with a scalar quantizer, under a fixed overload probability. Their optimization problem can be solved numerically, but no closed form solution was given. In [8] and [15] the design of an optimal *unconstrained* sigma-delta filter was studied, under the assumption of a fixed scalar quantizer which can only be scaled in order to control the overload probability. Equations that characterize the optimal filter were derived. However, the obtained expressions usually yield filters with an infinite number of taps, and do not provide the solution to the constrained problem. It is also worth mentioning that for the case of a stationary Gaussian process  $\{X_n\}$  with  $L = 1$  (sampling at the Nyquist rate) and *known PSD*, the optimal infinite length filter under the assumption of high-resolution quantization is known to equal the optimal prediction filter of  $X_n$  from its (clean) past [12]. As already mentioned in the introduction, the high-resolution assumption never holds when  $L > 1$  and therefore this result is inapplicable for over-sampled signals.

*Proof of Theorem 2:* By Proposition 1, if the test channel

from Figure 1 achieves MSE distortion  $D$ , we must have

$$\sigma_{\Sigma\Delta}^2 = \frac{D}{\frac{1}{2\pi} \int_{-\pi/L}^{\pi/L} |1 - C(\omega)|^2 d\omega}.$$

By Theorem 1, the corresponding mutual information  $I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta})$  is equal to the mutual information  $I(U_n^{\text{DPCM}}; U_n^{\text{DPCM}} + N_n^{\text{DPCM}})$  in the DPCM test channel from Figure 2 with  $X_n^{\text{DPCM}} = S_n$ ,  $N_n^{\text{DPCM}} = W_n$  and  $\sigma_{\text{DPCM}}^2 = L \cdot D$ . Thus,

$$\begin{aligned} I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) &= I(U_n^{\text{DPCM}}; U_n^{\text{DPCM}} + N_n^{\text{DPCM}}) \\ &= \frac{1}{2} \log \left( 1 + \frac{\mathbb{E}(S_n - c_n * (S_n + W_n))^2}{L \cdot D} \right), \end{aligned} \quad (15)$$

where we have used (8), (10), and (12), to arrive at (15). It follows that among all filters in  $\mathcal{C}$ , the filter that minimizes (15) is  $C_D^*(Z)$ , and that it attains (14) with equality. ■

It is interesting to note [9] that since  $\{W_n\}$  is an i.i.d. process with variance  $L \cdot D$  and  $C(Z)$  is strictly causal, the mutual information (15) can also be written as

$$\begin{aligned} I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) &= \frac{1}{2} \log \left( \frac{\mathbb{E}(S_n + W_n - c_n * (S_n + W_n))^2}{L \cdot D} \right). \end{aligned} \quad (16)$$

Thus, the optimal predictor of  $S_n$  from the past of  $\{S_n + W_n\}$  is identical to the optimal predictor of  $S_n + W_n$  from its past samples. When  $C(Z)$  is taken as the (unique) infinite order optimal one-step prediction filter of  $S_n + W_n$  from its past samples, the prediction error variance is the entropy power of the process  $\{S_n + W_n\}$  [16], which equals

$$\frac{1}{2\pi} \int_{-\pi}^{\pi} \log(S_S(\omega) + L \cdot D) d\omega = (L \cdot D) \left( 1 + \frac{\sigma_X^2}{D} \right)^{1/L}. \quad (17)$$

Moreover, the infinite order prediction error

$$E_n^{\text{pred}} \triangleq S_n + W_n - c_n * (S_n + W_n)$$

is in this case a white process. This, together with (17) implies that for the optimal unconstrained sigma-delta filter  $C(Z)$  we must have

$$\begin{aligned} S_{E^{\text{pred}}}(\omega) &\triangleq |1 - C(\omega)|^2 (L \cdot D + S_S(\omega)) \\ &= (L \cdot D) \left( 1 + \frac{\sigma_X^2}{D} \right)^{1/L}, \quad \forall \omega \in [-\pi, \pi] \end{aligned} \quad (18)$$

Combining (16), (17), and (18) yields the following corollary.

*Corollary 1:* Let  $\{X_n^{\Sigma\Delta}\}$  be a Gaussian stationary process with variance  $\sigma_X^2$  whose PSD is zero for all  $\omega \notin [-\pi/L, \pi/L]$ . If the test channel from Figure 1 attains MSE distortion  $D$ , then

$$I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) \geq \frac{1}{2L} \log \left( 1 + \frac{\sigma_X^2}{D} \right). \quad (19)$$

with equality if and only if  $C(Z)$  is a strictly causal filter satisfying

$$|1 - C(\omega)|^2 = \begin{cases} \left( 1 + \frac{\sigma_X^2}{D} \right)^{-(L-1)/L} & \omega \in [-\frac{\pi}{L}, \frac{\pi}{L}] \\ \left( 1 + \frac{\sigma_X^2}{D} \right)^{1/L} & \omega \notin [-\frac{\pi}{L}, \frac{\pi}{L}], \end{cases} \quad (20)$$

and

$$\sigma_{\Sigma\Delta}^2 = \frac{D}{\frac{1}{2\pi} \int_{-\pi/L}^{\pi/L} |1 - C(\omega)|^2 d\omega} = \frac{L \cdot D}{\left( 1 + \frac{\sigma_X^2}{D} \right)^{-(L-1)/L}}.$$

*Remark 2:* Note that the existence of a strictly causal filter  $C(Z)$  which satisfies (20) is guaranteed by Wiener's spectral-factorization theory [16] due to the readily verified fact that

$$2^{\frac{1}{2\pi} \int_{-\pi}^{\pi} \log |1 - C(\omega)|^2 d\omega} = 1.$$

The optimal filter induces a two-level frequency response for  $|1 - C(\omega)|^2$ . In [11] Østergaard and Zamir used sigma-delta modulation to attain the optimal multiple-description rate-distortion region. Interestingly, the optimal filter  $C(Z)$  in their scheme also induced a two-level response for  $|1 - C(\omega)|^2$ . We also note that the optimality of the unconstrained filter specified by (20) can be deduced as a special case of [8, Section IV].

*Remark 3:* Note that for the optimal unconstrained filter  $C(Z)$  specified by (20), the pre- and post-filters from Figure 3 have no effect as long as the PSD of the input signal  $\{X_n^{\Sigma\Delta}\}$  is zero for all  $\omega \notin [-\pi/L, \pi/L]$ . However, filters with a finite number of taps will never incur a flat frequency response in the interval  $[-\pi/L, \pi/L]$ , and for such filters the systems from Figure 1 and Figure 3 will not be equivalent.

*Remark 4:* The output of the test channel from Figure 1 (as well as that from Figure 2) is of the form  $\hat{X}_n^{\Sigma\Delta} = X_n^{\Sigma\Delta} + E_n^{\Sigma\Delta}$ , where  $E_n^{\Sigma\Delta}$  has zero mean and variance  $D$ , and is statistically independent of  $X_n^{\Sigma\Delta}$ . This estimate can be further improved by applying scalar MMSE estimation for  $X_n^{\Sigma\Delta}$  from  $\hat{X}_n^{\Sigma\Delta}$ . This boils down to producing the estimate  $\hat{\hat{X}}_n^{\Sigma\Delta} = \alpha \hat{X}_n^{\Sigma\Delta}$ , where

$$\alpha = \frac{\sigma_X^2}{\sigma_X^2 + D}.$$

Consequently, the obtained MSE distortion is reduced to

$$\tilde{D} = \mathbb{E}(X_n^{\Sigma\Delta} - \alpha \hat{X}_n^{\Sigma\Delta})^2 = \frac{\sigma_X^2 \cdot D}{\sigma_X^2 + D}.$$

It is straightforward to verify [17] that with this improvement, the sigma-delta test channel from Figure 1 with  $C(Z)$  and  $\sigma_{\Sigma\Delta}^2$  as specified in Corollary 1 attains

$$I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + V_n^{\Sigma\Delta}) = \frac{1}{2L} \log \left( \frac{\sigma_X^2}{D} \right),$$

which is the optimal rate-distortion function for a stationary Gaussian source  $\{X_n^{\Sigma\Delta}\}$  with PSD as in (7). It follows that the sigma-delta test channel from Figure 1 with  $C(Z)$  and  $\sigma_{\Sigma\Delta}^2$  as specified in Corollary 1 is minimax optimal for the class of all stationary Gaussian sources with variance  $\sigma_X^2$  and PSD that

equals zero for all  $\omega \notin [-\pi/L, \pi/L]$ , i.e., no other system can achieve MSE distortion  $\tilde{D}$  with a smaller compression rate, universally for all sources in this class.

#### A. Extension to Frequency-Weighted Mean Squared Error Distortion

In many applications, higher values of distortion are acceptable in certain frequency bands while smaller distortion is permitted in other bands. The MSE distortion measure is inadequate for such scenarios, and a commonly used distortion measure, that (partially) captures such perceptual effects, is the frequency-weighted mean squared error (FWMSE) criterion. Under this criterion, the distortion is measured as

$$D_{\text{FWMSE}} \triangleq \frac{1}{2\pi} \int_{-\pi}^{\pi} P(\omega) S_E(\omega) d\omega, \quad (21)$$

where  $P(\omega)$  is a non-negative weight function, and  $S_E(\omega)$  is the PSD of the error process  $E_n \triangleq X_n^{\Sigma\Delta} - \hat{X}_n^{\Sigma\Delta}$ . Note that for  $P(\omega) = 1, \forall \omega \in [-\pi, \pi]$ , the FWMSE criterion reduces to the MSE one. The next theorem shows that the constrained optimal sigma-delta filter under the FWMSE criterion is the optimal constrained prediction filter of a noisy process defined according to the weight function  $P(\omega)$ .

**Theorem 3:** Let  $\{X_n^{\Sigma\Delta}\}$  be a Gaussian stationary process with variance  $\sigma_X^2$  whose PSD is zero for all  $\omega \notin [-\pi/L, \pi/L]$ ,  $P(\omega)$  a weighting function which forms a valid PSD, and  $\mathcal{C}$  a family of strictly causal filters. Define the “virtual” process  $\{S_n\}$  as a Gaussian stationary process with PSD

$$S_X^{\text{FWMSE}}(\omega) = \begin{cases} L\sigma_X^2 P(\omega) & \text{for } |\omega| \leq \pi/L \\ 0 & \text{for } \pi/L < |\omega| < \pi \end{cases}, \quad (22)$$

and the “virtual” process  $\{W_n\}$  as a Gaussian i.i.d. random process statistically independent of  $\{S_n\}$  with variance  $L \cdot D_{\text{FWMSE}}$ ,  $D_{\text{FWMSE}} > 0$ . Let

$$\sigma_{D_{\text{FWMSE}}}^{*2} = \min_{C(Z) \in \mathcal{C}} \mathbb{E} (S_n - c_n * (S_n + W_n))^2$$

$$C_{D_{\text{FWMSE}}}^*(Z) = \operatorname{argmin}_{C(Z) \in \mathcal{C}} \mathbb{E} (S_n - c_n * (S_n + W_n))^2.$$

If the filter  $C(Z)$  in the sigma-delta test channel from Figure 1 belongs to  $\mathcal{C}$  and the FWMSE distortion w.r.t.  $P(\omega)$  attained by this test channel is  $D_{\text{FWMSE}}$ , then

$$I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) \geq \frac{1}{2} \log \left( 1 + \frac{\sigma_{D_{\text{FWMSE}}}^{*2}}{L \cdot D_{\text{FWMSE}}} \right),$$

with equality if  $C(Z) = C_{D_{\text{FWMSE}}}^*(Z)$ .

*Sketch of proof:* The proof is fairly similar to that of Theorem 2. Thus, for brevity, we omit the full proof and only highlight its main steps:

- Repeat the derivation of Proposition 1 where now the MSE distortion is replaced by FWMSE distortion. Note that this has no effect on  $I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta})$ .
- Repeat the derivation of Proposition 2 where the PSD of the input process is (22), rather than (7). Note that this changes  $I(U_n^{\text{DPCM}}; U_n^{\text{DPCM}} + N_n^{\text{DPCM}})$ , but has no effect on the attained distortion.

- It follows that the DPCM test channel for the process  $\{S_n\}$  under MSE distortion is equivalent to the sigma-delta test channel with input  $\{X_n^{\Sigma\Delta}\}$  under FWMSE distortion, in the sense that in both channels if the attained distortion is  $D_{\text{FWMSE}}$  (under the appropriate distortion measure), then

$$I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) = I(U_n^{\text{DPCM}}; U_n^{\text{DPCM}} + N_n^{\text{DPCM}})$$

$$= \frac{1}{2} \log \left( 1 + \frac{\mathbb{E} (S_n - c_n * (S_n + W_n))^2}{L \cdot D_{\text{FWMSE}}} \right).$$

#### B. Sigma-Delta Modulation with an Interleaved Vector Quantizer

The goal of this short subsection is to give the test channel from Figure 1 an operational meaning, i.e., to show how the AWGN from the figure can be replaced with a lossy source code of rate  $R = I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta})$  whose incurred quantization noise is distributed as  $N_n^{\Sigma\Delta}$ . As already mentioned, the key idea is to use an interleaver [9]–[11], as we now recall.

Assume that  $\{X_n^{\Sigma\Delta}\}$ , the input process to the sigma-delta modulator, has a decaying memory, such that  $X_n^{\Sigma\Delta}$  is essentially independent of all samples of sufficiently distant sampling times. In order to compress an  $N$ -dimensional vector

$$\mathbf{x}^{\Sigma\Delta} = [X_1^{\Sigma\Delta}, \dots, X_N^{\Sigma\Delta}],$$

containing  $N$  consecutive samples of the process  $\{X_n^{\Sigma\Delta}\}$ , we first split it into  $K$  vectors

$$\mathbf{x}_k^{\Sigma\Delta} = [X_{(k-1)M+1}^{\Sigma\Delta}, \dots, X_{kM}^{\Sigma\Delta}], \quad k = 1, \dots, K,$$

where  $M \triangleq N/K$ . Now, we can apply  $K$  parallel sigma-delta modulators, one for each such vector, where the only coupling between the  $K$  parallel systems is through the quantization step, which is applied jointly on all of them, as depicted in Figure 4. By our assumption that  $\{X_n^{\Sigma\Delta}\}$  has decaying memory, if  $M$  is large enough the  $K$  inputs that enter the quantizer  $\underline{Q}(\cdot) = [Q_1(\cdot), \dots, Q_K(\cdot)]$  are i.i.d. random variables distributed as  $U_n^{\Sigma\Delta}$  from Figure 1. For large enough  $K$ , standard rate-distortion arguments imply that there exists a vector quantizer with rate  $I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta})$  that induces quantization noise distributed as  $N_n^{\Sigma\Delta}$ .

### III. SIGMA-DELTA MODULATION WITH A SCALAR UNIFORM QUANTIZER

The previous subsection showed how to replace the AWGN channel in Figure 1 with a vector quantizer whose rate is arbitrarily close to  $R = I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta})$  and whose induced quantization noise is distributed as  $N_n^{\Sigma\Delta}$ . The inputs to the vector quantizer are vectors of i.i.d. Gaussian components. Thus, any “off-the-shelf” rate-distortion optimal vector quantizer for an i.i.d. Gaussian source can be used. The total sigma-delta compression system that is obtained is therefore simple in the sense that it only requires the vector quantizer to be good for quantizing an i.i.d. Gaussian source, which is



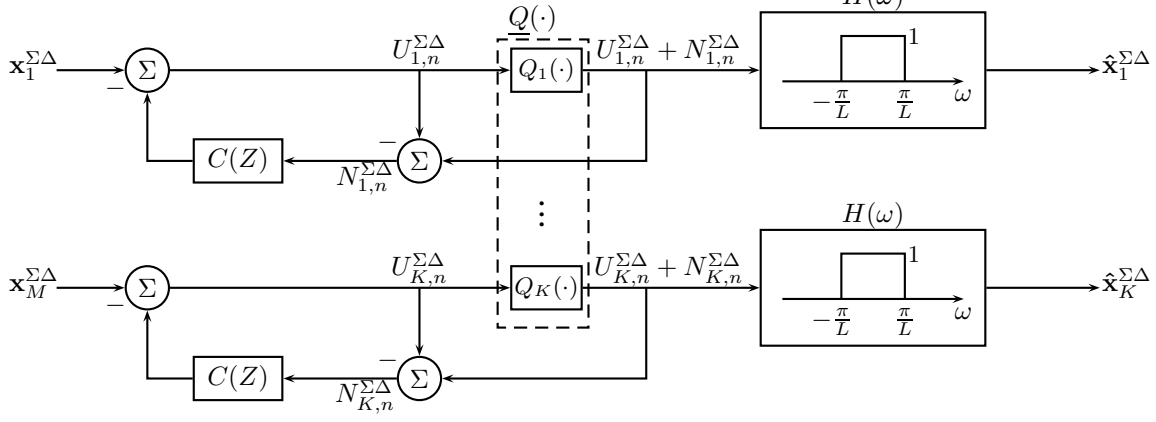


Fig. 4.  $K$  parallel sigma-delta modulators coupled by an  $K$ -dimensional quantizer  $\underline{Q}(\cdot)$ .

a standard task, rather than requiring it to be a good quantizer for a band-limited Gaussian source.

However, the sigma-delta modulation architecture is mainly used for A/D and D/A conversion. In such applications, vector quantization is typically out of the question, and simple uniform scalar quantizers of finite support are used. For such quantizers, the quantization error is composed of two main factors [1]: *granular errors* that correspond to the quantization error in the case where the input signal falls within the quantizer's support, and *overload errors* that correspond to the case where the input signal falls outside the quantizer's support. Due to the feedback loop, inherent to the sigma-delta modulator, errors of the latter kind, whose magnitude is not bounded, may have a disastrous effect as they jeopardize the system's stability. In order to avoid such errors, the support of the quantizer has to be chosen appropriately. As the support of the quantizer determines its rate for a given quantization resolution, the overload probability can be controlled by increasing the quantization rate.<sup>2</sup>

We shall show that, given that overload errors did not occur, the quantization noise can be modeled as an additive noise. Thus, the test channel from Figure 1 accurately predicts the total distortion incurred by a sigma-delta A/D (or D/A) in this case. Moreover, the overload probability is a doubly exponentially decreasing function of  $R - I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta})$ , where  $2^R$  are the number of levels in the scalar quantizer. Thus, fixing the desired overload error probability as  $P_{ol}$ , we may achieve the MSE distortion predicted by the test channel from Figure 1 (characterized in Proposition 1) with a scalar quantizer whose rate is  $I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) + \delta(P_{ol})$ , where  $\delta(P_{ol}) = \mathcal{O}\left(\log \log \left(\frac{1}{P_{ol}}\right)\right)$ .

Let  $Q_{R,\sigma^2}(\cdot)$  be a uniform quantizer with quantization step  $\sqrt{12\sigma^2}$  and  $2^R$  quantization levels, such that the quantizer support is  $[-\Gamma/2, \Gamma/2)$ , where  $\Gamma \triangleq 2^R \sqrt{12\sigma^2}$ , see Figure 5. Our goal is to analyze the distortion and overload probability attained by a sigma-delta modulator that uses a  $Q_{R,\sigma^2}(\cdot)$  quantizer, as a function of  $R$  and  $\sigma_{\Sigma\Delta}^2$ .

<sup>2</sup>As discussed in Section I-B, one can try to limit the effect of overload errors by placing various constraints on  $C(Z)$ . Here, we restrict attention to controlling the overload probability.

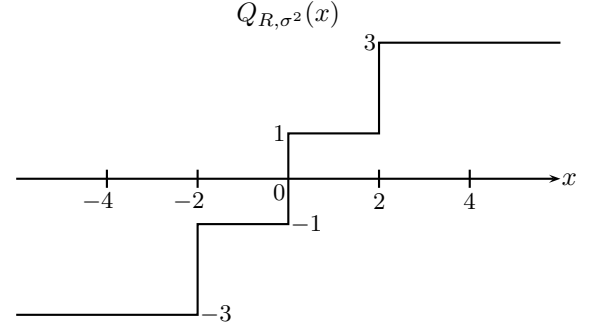


Fig. 5. An illustration of  $Q_{R,\sigma^2}(\cdot)$  for  $R = 2$  and  $\sigma^2 = 1/3$ .

Clearly, if we employ the scalar sigma-delta modulator on a long enough input sequence, an overload event will eventually occur. As discussed above, the effects of overload errors can be amplified due to the feedback loop, and in this case the average MSE may significantly grow. We therefore split the input sequence into finite blocks of length  $N$ , and initialize the memory of the filter  $C(Z)$  with zeros before the beginning of each new block. This makes sure that the effect of an overload error in the original system is restricted to the block where it occurs.

The analysis is made much simpler by introducing a subtractive *dither* [17]. Namely, let  $\{Z_n\}$  be a sequence of i.i.d. random variables uniformly distributed over the interval  $[-\sqrt{12\sigma_{\Sigma\Delta}^2}/2, \sqrt{12\sigma_{\Sigma\Delta}^2}/2)$ . In order to quantize  $U_n^{\Sigma\Delta}$ , we add  $Z_n$  to it before applying the quantizer, and subtract  $Z_n$  afterwards, such that the obtained result is  $Q_{R,\sigma_{\Sigma\Delta}^2}(U_n^{\Sigma\Delta} + Z_n) - Z_n$ . Adding and subtracting  $U_n^{\Sigma\Delta}$ , we get  $U_n^{\Sigma\Delta} + (Q_{R,\sigma_{\Sigma\Delta}^2}(U_n^{\Sigma\Delta} + Z_n) - (U_n^{\Sigma\Delta} + Z_n))$ , and the quantization error is therefore

$$N_n \triangleq Q_{R,\sigma_{\Sigma\Delta}^2}(U_n^{\Sigma\Delta} + Z_n) - (U_n^{\Sigma\Delta} + Z_n) \quad (23)$$

The main result in this section is the following.

**Theorem 4:** Let  $D$  be the MSE distortion attained by the test channel in Figure 1 with a filter  $C(Z)$  of finite length, and  $I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta})$  the scalar mutual information between the input and output of the AWGN channel in the same figure.

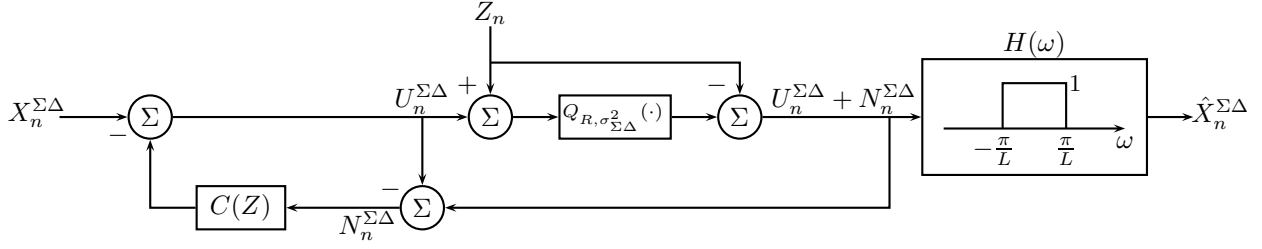


Fig. 6. A sigma-delta modulator with a dithered scalar uniform quantizer. The input is assumed to be over-sampled at  $L$  times the Nyquist rate, and the dither sequence  $\{Z_n\}$  is assumed to be an i.i.d. sequence of random variables uniformly distributed over the interval  $[-\sqrt{12\sigma_{\Sigma\Delta}^2}/2, \sqrt{12\sigma_{\Sigma\Delta}^2}/2]$  and statistically independent  $\{X_n^{\Sigma\Delta}\}$ .

For any  $0 < P_{ol} < 1$  the scalar sigma-delta modulator from Figure 6 applied on a sequence of  $N$  consecutive source samples with quantization rate  $R = I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) + \delta(P_{ol})$  attains MSE distortion smaller than

$$\frac{D(1 + o_N(1))}{1 - P_{ol}},$$

given that overload did not occur. In addition, the overload probability is smaller than  $P_{ol}$ , where  $o_N(1) \rightarrow 0$  as  $N$  increases, and

$$\delta(P_{ol}) \triangleq \frac{1}{2} \log \left( -\frac{2}{3} \ln \frac{P_{ol}}{2N} \right). \quad (24)$$

*Proof:* Let  $\tilde{Q}_{\sqrt{12\sigma^2}\mathbb{Z}}(x)$  be the operation of rounding  $x$  to the nearest point in the (infinite) lattice  $\sqrt{12\sigma^2}\mathbb{Z}$ . It is easy to verify that for any  $x \in [-\Gamma/2, \Gamma/2)$  we have

$$Q_{R, \sigma^2}(x) = \tilde{Q}_{\sqrt{12\sigma^2}\mathbb{Z}} \left( x + \frac{\sqrt{12\sigma^2}}{2} \right) - \frac{\sqrt{12\sigma^2}}{2}. \quad (25)$$

Applying (23) therefore yields that if overload did not occur in the  $n$ th sample, i.e., if  $|U_n^{\Sigma\Delta} + Z_n| \leq \Gamma/2$ , we have

$$N_n = \tilde{Q}_{\sqrt{12\sigma_{\Sigma\Delta}^2}\mathbb{Z}} \left( U_n^{\Sigma\Delta} + Z_n + \sqrt{3\sigma_{\Sigma\Delta}^2} \right) - \left( U_n^{\Sigma\Delta} + Z_n + \sqrt{3\sigma_{\Sigma\Delta}^2} \right). \quad (26)$$

Dealing with the overload event of the quantizer directly is rather involved. Instead, as done in [18], we first consider a *reference system* with an infinite-support quantizer ( $R = \infty$ ) and analyze its performance. If the magnitude of the input to the infinite-support quantizer never exceeds  $\Gamma/2$  within the processed block, then clearly the reference system is completely equivalent to the original system within this block. Thus, it suffices to find the average distortion of the reference system and the probability that the input to its quantizer exceeds  $\Gamma/2$  within a block. In what follows we will therefore assume that the quantization noise is given by (26) regardless of whether or not  $|U_n^{\Sigma\Delta} + Z_n| \leq \Gamma/2$ , and account for the overload probability later.

Assuming that the dither sequence  $\{Z_n\}$  is drawn statistically independent of the process  $\{X_n^{\Sigma\Delta}\}$ , the Crypto Lemma, see, e.g. [17, Lemma 4.1.1], implies that  $\{N_n\}$  is an i.i.d. sequence of random variables uniformly distributed over

the interval  $[-\sqrt{12\sigma_{\Sigma\Delta}^2}/2, \sqrt{12\sigma_{\Sigma\Delta}^2}/2]$ , statistically independent of  $\{X_n^{\Sigma\Delta}\}$ . Note that  $N_n$  has zero mean and variance  $\sigma_{\Sigma\Delta}^2$ . Following this reasoning, the reference sigma-delta data converter depicted in Figure 6 (with an infinite-support quantizer) is equivalent to the test channel from Figure 1 with  $N_n^{\Sigma\Delta} \sim \text{Uniform}([- \sqrt{12\sigma_{\Sigma\Delta}^2}/2, \sqrt{12\sigma_{\Sigma\Delta}^2}/2])$  instead of  $N_n^{\Sigma\Delta} \sim \mathcal{N}(0, \sigma_{\Sigma\Delta}^2)$ . Thus, the average MSE distortion attained by the reference scalar sigma-delta modulator from Figure 6 is as given in Proposition 1 up to a multiplicative factor of  $1 + o_N(1)$  that accounts for edge effects. These effects are the by-product of the operation of nulling the filter memory at the beginning of each new block, which incurs temporal non-stationarities. In particular, if the filter  $C(Z)$  has  $L$  taps, then only after  $L$  samples within the block the statistics of the process  $\{U_n^{\Sigma\Delta}\}$  will converge to its stationary distribution. However, if the block length is sufficiently large w.r.t. the filter length and the inverse of the MSE distortion, the influence of these effects vanishes.

Next, we turn to analyze the probability that an overload error occurs within a block of length  $N$ , as a function of  $R$  and  $I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta})$ . Since this event is equivalent to the event that at the reference system some input to the quantizer exceeds  $\Gamma/2$  in magnitude within the block, it suffices to upper bound the probability of the latter event.

Assume the reference scalar sigma-delta modulator from Figure 6 is applied to a vector  $\mathbf{x}^{\Sigma\Delta} = [X_1^{\Sigma\Delta}, \dots, X_N^{\Sigma\Delta}]$  of  $N$  consecutive samples of the process  $\{X_n^{\Sigma\Delta}\}$ , where the memory of the filter  $C(Z)$  is initialized with zeros. Define the event  $\text{OL}_k \triangleq \{|U_k^{\Sigma\Delta} + N_k^{\Sigma\Delta}| > \Gamma/2\}$  and the event  $\text{OL} \triangleq \cup_{k=1}^N \text{OL}_k$ . By the union bound, we have

$$P_{ol} \triangleq \Pr(\text{OL}) \leq \sum_{k=1}^N \Pr(\text{OL}_k). \quad (27)$$

The random variable  $U_k^{\Sigma\Delta} + N_k^{\Sigma\Delta} = X_k^{\Sigma\Delta} + (\delta_k - c_k) * N_k^{\Sigma\Delta}$  is a linear combination of a Gaussian random variable  $X_k^{\Sigma\Delta}$  and statistically independent uniform random variables  $\{N_k^{\Sigma\Delta}\}$ . In [19, Lemma 4] the probability that a random variable of this type exceeds a certain threshold was bounded in terms of its variance. Applying this bound to  $U_k^{\Sigma\Delta} + N_k^{\Sigma\Delta}$

yields

$$\begin{aligned} \Pr\left(|U_k^{\Sigma\Delta} + N_k^{\Sigma\Delta}| > \Gamma/2\right) &\leq 2 \exp\left\{-\frac{\Gamma^2}{8\mathbb{E}(U_k^{\Sigma\Delta} + N_k^{\Sigma\Delta})^2}\right\} \\ &= 2 \exp\left\{-\frac{12\sigma_{\Sigma\Delta}^2 2^{2R}}{8(\mathbb{E}(U_k^{\Sigma\Delta})^2 + \mathbb{E}(N_k^{\Sigma\Delta})^2)}\right\}, \end{aligned}$$

where in the last equality we have used the definition of  $\Gamma$  and the fact that  $U_k^{\Sigma\Delta}$  and  $N_k^{\Sigma\Delta}$  are statistically independent. Equivalently, we may write

$$\begin{aligned} \Pr\left(\text{OL}_k\right) &\leq 2 \exp\left\{-\frac{12\sigma_{\Sigma\Delta}^2 2^{2R}}{8\sigma_{\Sigma\Delta}^2 \left(1 + \frac{\mathbb{E}(U_k^{\Sigma\Delta})^2}{\sigma_{\Sigma\Delta}^2}\right)}\right\} \\ &= 2 \exp\left\{-\frac{3}{2} 2^{2\left(R - \frac{1}{2} \log\left(1 + \frac{\mathbb{E}(U_k^{\Sigma\Delta})^2}{\sigma_{\Sigma\Delta}^2}\right)\right)}\right\} \\ &= 2 \exp\left\{-\frac{3}{2} 2^{2(R - I(U_k^{\Sigma\Delta}; U_k^{\Sigma\Delta} + N_k^{\Sigma\Delta}))}\right\}, \quad (28) \end{aligned}$$

where we have used (5) in the last equality. Substituting (28) into (27) gives

$$P_{ol} \leq 2 \sum_{k=1}^N \exp\left\{-\frac{3}{2} 2^{2(R - I(U_k^{\Sigma\Delta}; U_k^{\Sigma\Delta} + N_k^{\Sigma\Delta}))}\right\}. \quad (29)$$

Note that  $\mathbb{E}(U_k^{\Sigma\Delta})^2 = \sigma_X^2 + \sigma_{\Sigma\Delta}^2 \sum_{m=1}^k c_k^2$  is monotonically nondecreasing in  $k$  and is given by (6) for values of  $k$  that are greater than the length of the filter  $c_k$ . We can therefore further bound (29) as

$$P_{ol} \leq 2N \exp\left\{-\frac{3}{2} 2^{2(R - I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}))}\right\}, \quad (30)$$

where  $I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta})$  is as given in Proposition 1. To summarize, we have shown that the reference system achieves the same MSE distortion  $D$  as characterized by Proposition 1 up to a  $1 + o_N(1)$  multiplicative term, and that the probability that one of the quantizer input samples exceeds  $\Gamma/2$  in magnitude within a block of length  $N$ , is bounded by (30). For our original system whose quantizer has finite support of  $[-\Gamma/2, \Gamma/2)$ , this means that the overload probability is also upper bounded by the RHS of (30). Moreover, the average distortion it achieves if overload did not occur is the same as that of the reference system conditioned on the event that OL did not occur. Denote this conditioned expected distortion by  $D_{\overline{\text{OL}}}$  and the expected distortion conditioned on the event that OL did occur by  $D_{\text{OL}}$ . For the reference system, we have

$$D(1 + o_N(1)) = \Pr(\overline{\text{OL}})D_{\overline{\text{OL}}} + \Pr(\text{OL})D_{\text{OL}} \geq \Pr(\overline{\text{OL}})D_{\overline{\text{OL}}},$$

and therefore

$$D_{\overline{\text{OL}}} \leq \frac{D(1 + o_N(1))}{1 - P_{ol}}.$$

This shows that the scalar sigma-delta system from Figure 6, whose quantizer has limited support  $[-\Gamma/2, \Gamma/2)$ , with  $R = I(U_n^{\Sigma\Delta}; U_n^{\Sigma\Delta} + N_n^{\Sigma\Delta}) + \delta(P_{ol})$  achieves the same average MSE distortion as the test channel from Figure 1 up to a multiplicative factor of  $(1 + o_N(1))/(1 - P_{ol})$ , with block error

probability smaller than  $2N \exp\left\{-\frac{3}{2} 2^{2\delta}\right\}$ . Thus, Proposition 1 characterizes the rate-distortion tradeoff achieved by the scalar sigma-delta system up to the aforementioned factor and a constant rate penalty  $\delta(P_{ol})$ , that depends on the target overload error probability. To be more precise, for any  $0 < P_{ol} < 1$ , taking the rate penalty as in (24) guarantees that the overload error probability is smaller than  $P_{ol}$ . ■

#### ACKNOWLEDGEMENTS

We thank Jan Østergaard and Ram Zamir for their valuable comments on an earlier version of this manuscript.

#### REFERENCES

- [1] N. S. Jayant and P. Noll, *Digital coding of waveforms: principles and applications to speech and video*. Englewood Cliffs, NJ: Prentice-Hall, 1984.
- [2] R. A. McDonald, "Signal-to-noise and idle channel performance of differential pulse code modulation systems particular applications to voice signals," *Bell System Technical Journal*, vol. 45, no. 7, pp. 1123–1151, 1966.
- [3] R. Gray and D. Neuhoff, "Quantization," *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2325–2383, Oct 1998.
- [4] C. Cutler, "Differential quantization of communication signals," 1952, US Patent 2,605,361.
- [5] —, "Transmission systems employing quantization," 1960, US Patent 2,927,962 (filed 1954).
- [6] S. Tewksbury and R. Hallock, "Oversampled, linear predictive and noise-shaping coders of order  $n > 1$ ," *IEEE Transactions on Circuits and Systems*, vol. 25, no. 7, pp. 436–447, Jul 1978.
- [7] M. Derpich, "Optimal source coding with signal transfer function constraints," Ph.D. dissertation, University of Newcastle, 2009.
- [8] M. Derpich, E. Silva, D. Quevedo, and G. Goodwin, "On optimal perfect reconstruction feedback quantizers," *IEEE Transactions on Signal Processing*, vol. 56, no. 8, pp. 3871–3890, Aug 2008.
- [9] R. Zamir, Y. Kochman, and U. Erez, "Achieving the Gaussian rate-distortion function by prediction," *IEEE Transactions on Information Theory*, vol. 54, no. 7, pp. 3354–3364, July 2008.
- [10] T. Guess and M. K. Varanasi, "An information-theoretic framework for deriving canonical decision-feedback receivers in Gaussian channels," *IEEE Transactions on Information Theory*, vol. IT-51, pp. 173–187, Jan 2005.
- [11] J. Østergaard and R. Zamir, "Multiple-description coding by dithered delta-sigma quantization," *IEEE Transactions on Information Theory*, vol. 55, no. 10, pp. 4661–4675, Oct 2009.
- [12] P. Noll, "On predictive quantizing schemes," *The Bell System Technical Journal*, vol. 57, no. 5, pp. 1499–1532, May 1978.
- [13] M. A. Gerzon and P. G. Craven, "Optimal noise shaping and dither of digital signals," in *Audio Engineering Society Convention 87*, 1989.
- [14] H. Spang III and P. Schultheiss, "Reduction of quantizing noise by use of feedback," *IRE Transactions on Communications Systems*, vol. 10, no. 4, pp. 373–380, Dec 1962.
- [15] M. Derpich and J. Østergaard, "Improved upper bounds to the causal quadratic rate-distortion function for Gaussian stationary sources," *IEEE Transactions on Information Theory*, vol. 58, no. 5, pp. 3131–3152, May 2012.
- [16] T. Berger, *Rate distortion theory: A mathematical basis for data compression*. Englewood Cliffs, NJ: Prentice-Hall, 1971.
- [17] R. Zamir, *Lattice Coding for Signals and Networks*. Cambridge: Cambridge University Press, 2014.
- [18] A. Ben-Yishai and O. Shayevitz, "The Gaussian channel with noisy feedback: Near-capacity performance via simple interaction," 2014. [Online]. Available: <http://arxiv.org/abs/1407.8022>
- [19] O. Ordentlich and U. Erez, "Precoded integer-forcing universally achieves the MIMO capacity to within a constant gap," *IEEE Transactions on Information Theory*, vol. 61, no. 1, pp. 323–340, Jan 2015.